



# OCR-D

KOORDINIERUNGSPROJEKT ZUR  
WEITERENTWICKLUNG VON OCR-VERFAHREN

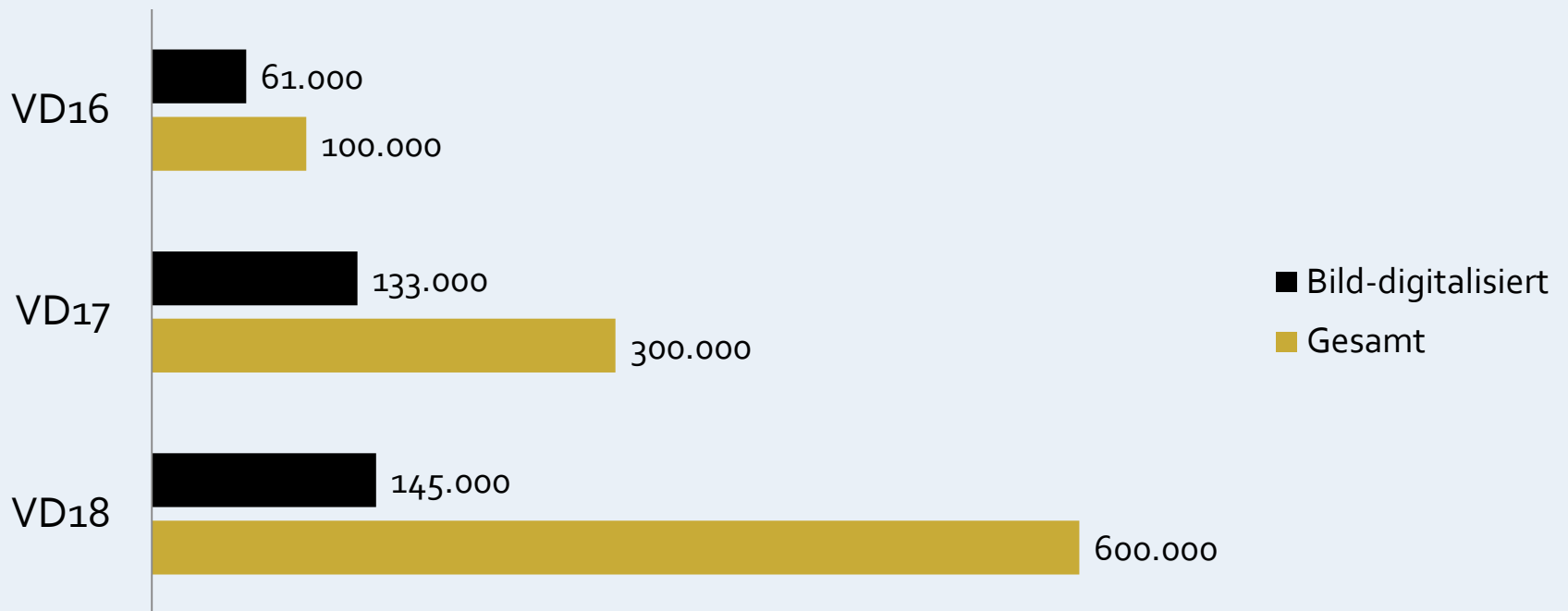
Gefördert von der Deutschen Forschungsgemeinschaft

25.02.2016

Elisa Herrmann



## VD 16-18



\* Gerundete Zahlen



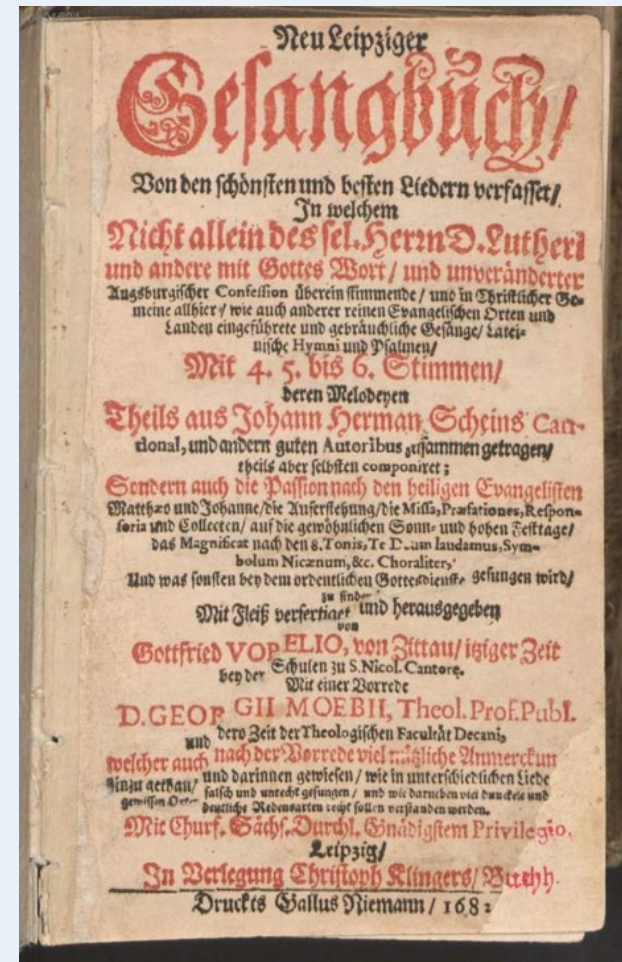
# Ziel

**Konzeptionelle Vorbereitung der Transformation der VD-  
Drucke (16.-18. Jh.) und der Drucke des 19. Jh. in  
maschinenlesbare Form.**



# Herausforderungen

- Material:
  - Sprachen: v.a. Latein, Deutsch
  - Schriftarten und -ausprägungen: u.a. Antiqua, Fraktur, Kursive
  - verschiedene Textsorten mit spezifischem Layout
- Uneinheitliche Standards
- Neuprozessierung, dynamische Datensicherung



Digitalisierung gefördert durch die Deutsche Forschungsgemeinschaft - DFG

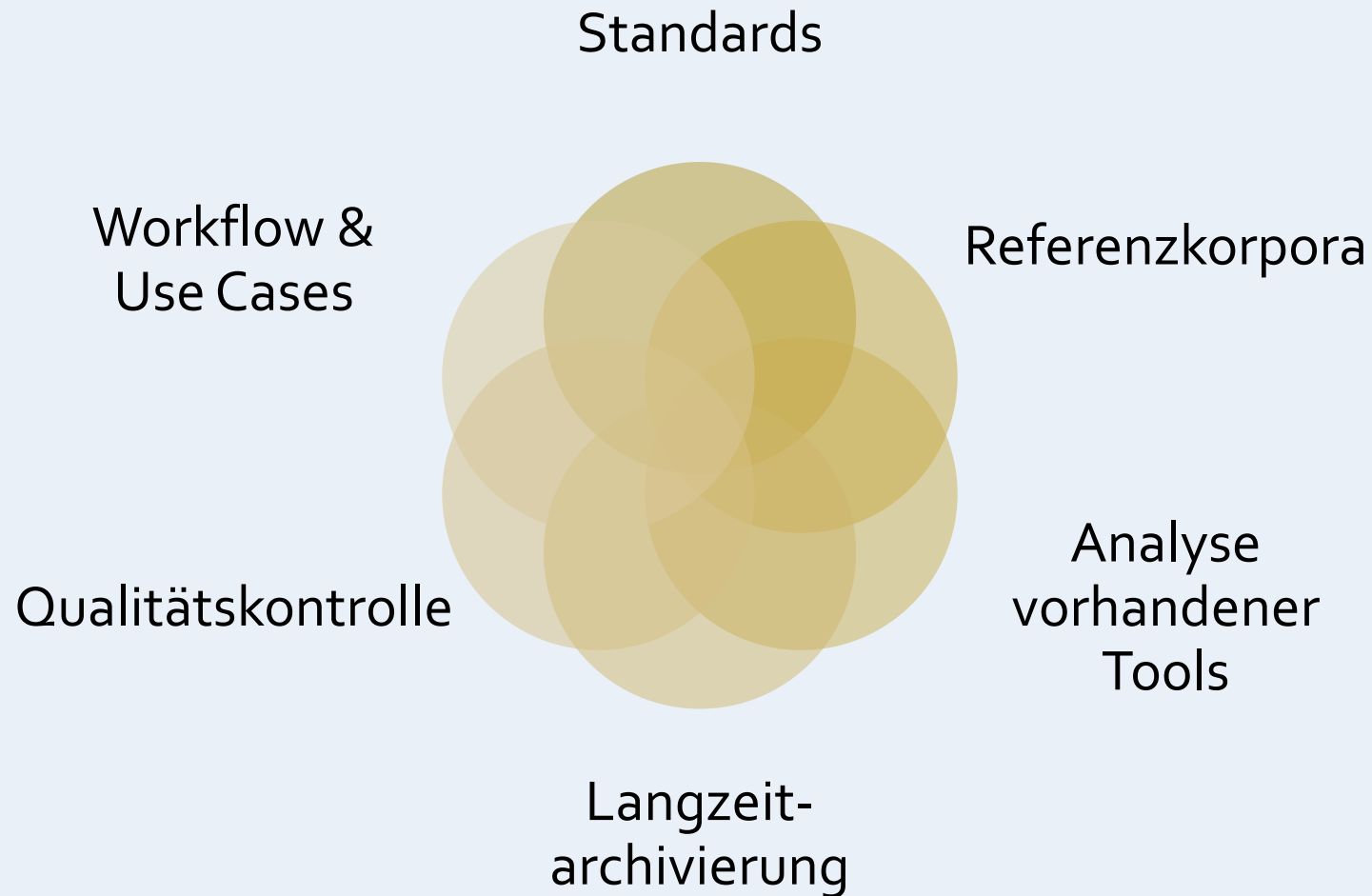


# Das Projekt

- Projektpartner
  - Herzog August Bibliothek Wolfenbüttel
  - Berlin-Brandenburgische Akademie der Wissenschaften, insb. Deutsches Textarchiv (DTA)
  - Bayerische Staatsbibliothek in München
- 2 Phasen:
  1. Aufbau der Koordinierungsstruktur und Konzeption der Projektphase
  2. Ausschreibung und konzeptionelle Begleitung der Pilotprojekte



# Arbeitspakete





# Funktionsmodell


- Einbindung neuer Erkenntnisse aus den einzelnen Arbeitspaketen
- Einbindung bestehender Werkzeuge
- Modularer Aufbau
  - Adaptive Anpassung an verschiedene Bedürfnisse
  - Anpassung an zukünftige Entwicklungen



# Funktionsmodell – I



- Bild-Digitalisat bereitstellen



- Splitting (optional)



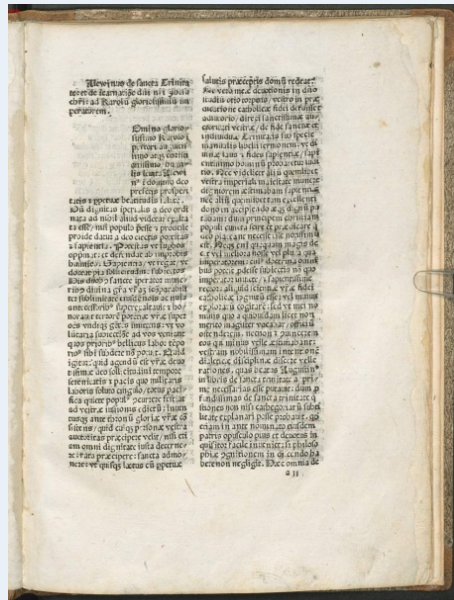
- Bildvorsortierung



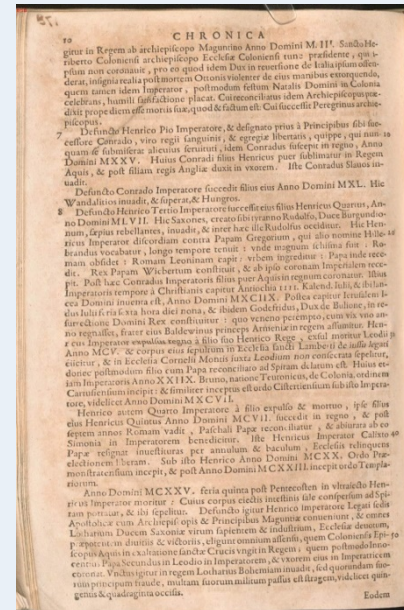


# Bildvorsortierung

- Bildvorsortierung
  - Z.B. Layoutanalyse



[http://daten.digital-sammlungen.de/bsb00005199/image\\_9](http://daten.digital-sammlungen.de/bsb00005199/image_9)



[http://daten.digital-sammlungen.de/bsb00010079/image\\_16](http://daten.digital-sammlungen.de/bsb00010079/image_16)



# Funktionsmodell – II

- 
- Preprocessing Page Level
  - (z.B. Cropping, Deskewing, Binarisierung)

- 
- Qualitätskontrolle

- 
- OLR1: Page Segmentation

- 
- Preprocessing Segment Level

- 
- Qualitätskontrolle



# Page Segmentation



- Segmentierung des Bilddigitalisats in
  - Textzonen
  - Nichttextzonen



# Funktionsmodell - III



- OCR



- Qualitätskontrolle



- OLR 2: Region Classification



- Qualitätskontrolle



# OCR & OLR

- OCR
  - Vgl. bzw. Kombination **klassischer Zeichenerkennungsverfahren** auf Glyphenebene mit **segmentierungsfreien Ansätzen** (z.B. "Deep Learning"-Verfahren auf Basis neuronaler Netze)
- OLR
  - **Region Classification:** Bestimmung der layout-semantischen Funktion der einzelnen Regionen (Überschrift, Marginalie, Fußnote etc.)
  - **Document Analysis:** Extrapolierung der Dokumentstruktur aus den entsprechenden Strukturelementen (Überschrift)



# Funktionsmodell - IV



- Export

- Nachkorrektur/ Crowdsourcing

- Qualitätskontrolle

- LZA



- Antworten auf technische, informationswissenschaftliche & organisatorische Herausforderungen

**Konsolidiertes Verfahren zur OCR-Verarbeitung von Digitalisaten des schriftlichen deutschen Kulturerbes des 16.-19. Jh.**



# Kontakt Daten

- Webseite: [www.ocr-d.de](http://www.ocr-d.de)



- Elisa Herrmann, [elisa.herrmann@hab.de](mailto:elisa.herrmann@hab.de), +49 5331 808-306